



# Pipeline Pilot Interface to *FTrees* Fragment Spaces

Version 2.3.0.1

## User Guide

(for *FTrees* version 2.3.0 and above and Pipeline Pilot version 6.1 and above)

Edgar Derksen, Sally Hindle



The idea of Feature Trees was born in 1997 during Matthias Rarey's six-month research stay at SmithKline Beecham Pharmaceuticals R&D, King of Prussia, PA (USA) and then further developed at Institute for Algorithms and Scientific Computing (SCAI), then part of the German National Research Center for Information Technology (GMD) and now the Fraunhofer Gesellschaft (FhG). Since 2002 BioSolveIT GmbH has been responsible for the licensing and continuing development of the *FTrees* software.

At this point we would like to thank Scott Dixon (SmithKline Beecham, now Metaphorics LCC.), Markus Wagener (SmithKline Beecham, now N.V. Organon), and Jens Lösel for a lot of helpful and constructive discussions during Matthias Rarey's stay at SmithKline Beecham and afterwards. Without them, the idea of Feature Trees would not have been evolved in the way it has. Also, Matthias Rarey thanks the GMD and SmithKline Beecham for funding his research stay in King of Prussia.

In summer 2000, the Feature Tree comparison algorithms were extended to search directly in large combinatorial chemistry spaces. A two-stage dynamic programming algorithm enables searching directly in chemistry spaces without an explicit enumeration of molecules [3]. This work was also done during a research stay in the US, this time at Roche Bioscience in Palo Alto. The chemistry space search algorithm was developed in cooperation with Martin Stahl (Hoffmann-La Roche, Basel) and we would like to thank him for this excellent cooperation. We also wish to thank Hans-Joachim Böhm, Hans Maag (both Roche), and Thomas Lengauer (GMD) for making this research stay possible.

Since then the Feature Trees software has been further developed and extended by several contributors including Marc Zimmermann (FhG - MTrees and the new Dynamic Matchsearch algorithm), Robert Fischer, Sally Hindle, and other developers at BioSolveIT GmbH and the Center for Bioinformatics (ZBH), University of Hamburg.

*This document contains proprietary information of BioSolveIT GmbH and is protected by copyright. It is provided together with Software of BioSolveIT under a license agreement and may be used only in accordance with the terms and conditions of this agreement. The document serves solely for the purpose of using the Software. No part of the document may be transferred to any third party or reproduced as a whole or in parts without written permission from BioSolveIT.*

*Base software: © 2001 by Fraunhofer Gesellschaft (FhI-SCAI); Getline library: © 1993 by Chris Thewalt; PVM library: © 1997 by University of Tennessee, Knoxville TN; Python library: © 1991-1995 by Stichting Mathematisch Centrum, Amsterdam, The Netherlands.*

# Contents

<b>Contents</b>	<b>3</b>
<b>1 Quick Start Steps</b>	<b>5</b>
1.1 Drag, drop and set license . . . . .	5
1.2 Note about the <i>FTrees</i> installation . . . . .	5
<b>2 Introduction</b>	<b>7</b>
2.1 About <i>FTrees</i> . . . . .	7
2.2 <i>FTrees-FS</i> Component in Pipeline Pilot . . . . .	7
2.3 The Example Fragment Space: <code>tfs</code> . . . . .	9
<b>3 The Internal and External <i>FTrees</i> Installations</b>	<b>11</b>
3.1 Using the Internal Installation . . . . .	11
3.2 Using an External <i>FTrees</i> Installation . . . . .	11
3.2.1 To Connect to an <i>FTrees</i> Installation on the Pipeline Pilot Server . . . . .	12
3.2.2 To Connect to an <i>FTrees</i> Installation on a Remote Linux Server . . . . .	13
<b>4 Trouble Shooting</b>	<b>15</b>
4.1 Problems connecting to the external installation . . . . .	15
4.2 Problems using the <i>ssh</i> Method . . . . .	16
4.3 Further help . . . . .	16
<b>5 Tips and Tricks</b>	<b>19</b>
5.1 Other Significant Parameters in the <i>FTrees</i> Components . . . . .	19
5.1.1 <Has Same File System> . . . . .	19
5.2 Accessing Other Domains within Pipeline Pilot . . . . .	19
<b>References</b>	<b>21</b>

**Bibliography**

# Quick Start Steps

## 1.1 Drag, drop and set license

1. Unzip the download package and save the contents on your system: the package contains the *FTrees* Fragment Spaces (*FTrees-FS*) in Pipeline Pilot component and a ready-to-use Fragment Space – do not unzip the Fragment Space file `knowledgespace-x.x.x-indep.zip`!
2. Drag and drop the component (the `.xml` file) onto the component area of your Pipeline Pilot client.
3. Add your license for the executable of *FTrees* in the field for the parameter `<Run FTrees -> on PP Server -> License Server or License File>` in the Implementation tab – if your license is available from a license server, simply type in the name of the server in this format `@servername`, or if you have a license file then you may browse for it using the `...` facility. *Note:* You must have a license for the Fragment Spaces module to use this component! See below and sections 3.2.1 and 3.2.2 for more details on the *FTrees* installation and license.
4. Enter the path to the example Fragment Space delivered with the component, for example `C:\My Documents\tfs_bundle.zip`
5. Enter the path to your query file (a molecule or Feature Tree file)
6. Use the component to generate molecules for a pipeline.

## 1.2 Note about the *FTrees* installation

The *FTrees* components are set by default to use a so-called internal installation of *FTrees*. This works as follows. When you run any of the *FTrees* components (and the parameter `<Run FTrees>` is set to *on PP Server* and the parameter `<Run FTrees-> on PP Server -> Use>` is set to *FTrees Auto Installation*), they search for a *FTrees* installation in the directory:

```
<scitegic install directory>/public/bin/BioSolveIT/
```

and use this installation to run the calculation using the license you gave in the third step described above. If a *FTrees* installation is not found in this directory, Pipeline Pilot will

install it itself using information from the component. Don't worry, you do not need Administrator rights to do this as it is actually Pipeline Pilot that carries out the installation and not you as a user.

Alternatively, you can also use an existing external installation of the *FTrees* software. This means you must already have *FTrees* installed on your system outside of Pipeline Pilot. If you have not already done so, visit the download page at BioSolveIT:

<http://www.biosolveit.de/download>

and fetch the download package for your system for the latest *FTrees* package. Follow the instructions in the package to install *FTrees* and receive your licenses.

See sections 3.2.1 and 3.2.2 for more details on the *FTrees* installation. **To use *FTrees* you must have a license for the Fragment Space module** [FFS].

# Introduction

## 2.1 About *FTrees*

*FTrees* is a piece of software for calculating the Feature Tree descriptor and comparing two or more of these descriptors to each other. The theory behind Feature Trees can be found in [2].

Rather than being based on a linear description such as bit strings or vectors, the Feature Tree descriptor represents the molecule as an unrooted tree where the nodes of the tree describe the major building blocks of the molecule. The comparison of two Feature Trees then proceeds using a recursive matching algorithm, splitting the trees into smaller and smaller subtrees. The Feature Tree approach has several advantages, the most important being the fact that the alignment of two Feature Trees can be translated into a comprehensible mapping of the two underlying molecules. For more details on the algorithms and achieved results see [2].

The very nature of the Feature Tree means it is also perfectly suited to working with fragments of molecules. Molecular fragments can be described by a Feature Tree just like a complete molecule and, therefore, can also be compared to each other. Together with the rules describing how fragments are allowed to be joined, the molecule fragments form a *FTrees* Fragment Space. Complete molecules can be used to search in the Fragment Space using the joining rules to guide the construction of fragments, see Figure 2.1 for an illustration. This technique is especially advantageous for searching in combinatorial sets of molecules like combinatorial libraries or combinatorial chemistry spaces. [3] contains more complete details about how the search algorithm works.

## 2.2 *FTrees-FS* Component in Pipeline Pilot

The *FTrees-FS* component can be thought of as a molecule generator for Pipeline Pilot. The component has no input port but takes a query file as a parameter. The results of searching with your queries is a set of resulting molecules output into the pipeline.

The component requires two files:

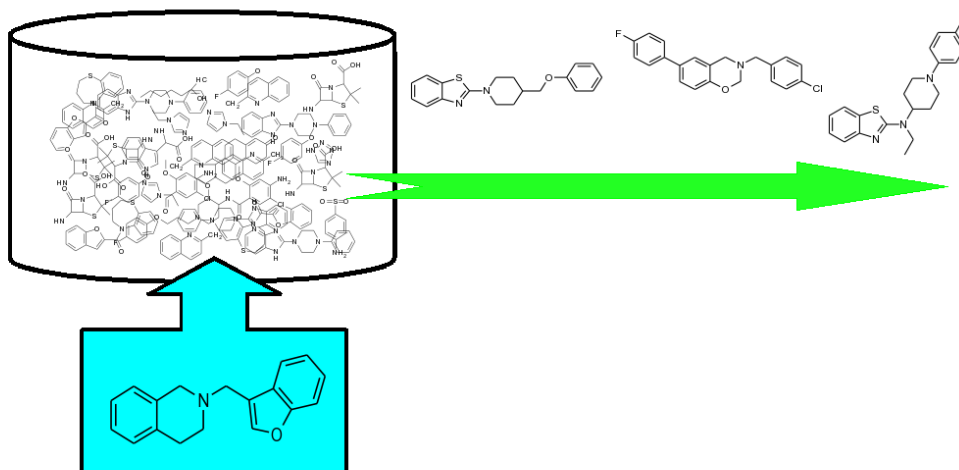


Figure 2.1: The *FTrees-FS* component takes a query molecule as input and uses it to generate new molecules from a Fragment Space based on a set of joining rules – the new molecules are sent via the Pass port into the pipeline

**Fragment Space File** A Fragment Space must be provided for the component as a `.zip` package. This is because a Fragment Space consists of several files and these *must* be kept together during any file transfer process. Inside the `.zip` package are the four files that form the Fragment Space in the format required by *FTrees-FS*. If you already have your own Fragment Spaces be sure to have all four files ready and pack them together using a zip program. The four files that form a Fragment Space in *FTrees-FS* format can be seen as an example in section 2.3.

**Query** The query file can be in any molecule input format readable by Pipeline Pilot, a Feature Tree file (`.fdf`) or the special Feature Tree database format used by the *FTrees* Writer and Reader components (`.ftdb`). You can use a multiple query file – make sure you adjust the maximum number of queries parameter.

The results of the search can be tuned according to how similar the resulting molecules are to the queries. Often, you do not want to see molecules almost exactly identical to your queries but are looking to generate new ideas. Change the parameter `<FTrees Options -> Target Similarity>` (an *FTrees* similarity value) to do this. If you find that the resulting molecules are too similar amongst each other, you can also set the diversity that must be seen between the resulting molecules. Change the parameter `<FTrees Options -> Target Diversity>`.

The resulting molecules exit the Pass port of the component. If you see any molecules exiting from the Fail port then these are queries that could not be interpreted correctly.

## 2.3 The Example Fragment Space: `tfs`

The *FTrees-FS* component is delivered with a ready-to-use Fragment Space `tfs`. This example Fragment Space made available by BioSolveIT contains approximately 60 000 molecule fragments from publically available data. The associated rules are based on the RECAP rules [1]. The size of the enumerated space is in the order of  $10^{13}$  for a maximum of three fragments and  $10^{23}$  for a maximum of five fragments!

The Fragment Space package `knowledgespace-x.x.x-indep.zip` must not be unpacked as the component will read the zip file directly.

Out of interest you may be curious about what is in the Fragment Space bundle. The following files are included:

`knowledgespace.fsf` This is the main Fragment Space file. It includes pointers to the molecule fragment file, Feature Tree fragment file, joining rule file and the joining rules themselves (linking rules)

`knowledgespace.mol2` The actual molecule fragments themselves. The fragments were generated using publically available data.

`knowledgespace.fdf` The Feature Tree descriptors for the molecule fragments.

`knowledgespace.ldf` The linkage rules for the fragments.



# The Internal and External *FTrees* Installations

## 3.1 Using the Internal Installation

As mentioned previously components are set by default to use a internal installation of *FTrees* (the parameter `<Run FTrees>` is set to *on PP Server* and the parameter `<Run FTrees-> on PP Server -> Use>` is set to *FTrees Auto Installation*). This is a default installation made whenever any *FTrees* components are used for the first time. The software is installed automatically by Pipeline Pilot within its own space on the server:

```
<scitegic install directory>/public/bin/BioSolveIT/
```

The installation is carried out once only (for that *FTrees* version) and only once per server not per user. The software is available to all users once it is installed. Note that a normal user cannot edit or delete this installation.

As this installation knows nothing about the licenses you may have for *FTrees*, you have to supply the license information separately. This is supplied using the parameter `<Run FTrees -> on PP Server -> License Server or License File>` as described above in the section 1.

## 3.2 Using an External *FTrees* Installation

You can also use an existing external installation of the *FTrees* software. This means you also have the opportunity to use settings different to those set by default. To do this, you must already have *FTrees* installed somewhere on your system outside of Pipeline Pilot. To install *FTrees* yourself, visit the download page at BioSolveIT:

```
http://www.biosolveit.de/download
```

and fetch the download package for your system for the latest *FTrees* package. Follow the instructions in the package to install *FTrees* and receive your licenses. Enter the license information for *FTrees* as described in the package **and not using the parameter `<Run FTrees -> on PP Server -> License Server or License File>` as for the internal installation.**

To use an external installation of *FTrees*, you must change the value of the parameter <Run FTrees-> on PP Server -> Use> in the Implementation tab to *preinstalled FTrees*.

There are actually two ways to use *FTrees* with an external installation. These are by using *FTrees* installed directly on the Pipeline Pilot server you are , or by accessing a remote machine where *FTrees* is installed using *ssh* method). The method is selected using the parameter <Run FTrees> Both methods are covered in more detail below.

### 3.2.1 To Connect to an *FTrees* Installation on the Pipeline Pilot Server

The most common scenario is that you will have an installation of *FTrees* on the Pipeline Pilot server. If you choose this option, you just have to enter the path to the executable and configuration file as parameters in the Implementation tab. Pipeline Pilot will then just start *FTrees* whenever it is required by making a call to the executable that you entered.

This method is given the name *on PP Server* – On Pipelinepilot Server. You can see an example in figure 3.1.

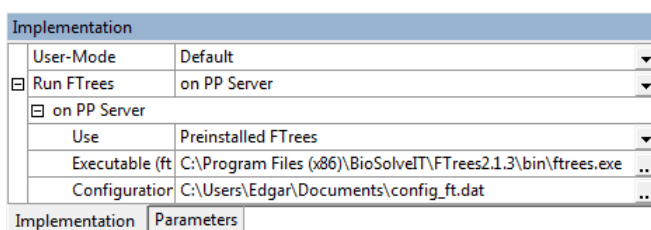


Figure 3.1: *ops*: Connect to an *FTrees* Installation on the Pipeline Pilot Server

- Requirement:
  - *FTrees* is installed on the Pipeline Pilot server. You can see which machine is the Pipeline Pilot server by starting your copy of Pipeline Pilot Client on your own workstation and find the name or IP of the server shown at the bottom right of the status bar (see figure 3.2). You must find where the *FTrees* installation is *on that machine*.
- General Steps:
  1. Set <Run FTrees> to *on PP Server*
  2. Set <Run FTrees-> Use> to *Preinstalled FTrees*
  3. Expand <opsFTreesExe>
  4. For the parameter <Run FTrees-> on PP Server -> Executable>, enter the path to the *FTrees* installation on Pipeline Pilot server. For example:  
C:\Programs\BioSolveIT\FTrees2\bin\ftrees.exe
  5. For the parameter <Run FTrees-> on PP Server -> Configuration>, enter the path to the *FTrees* configuration file config\_ft.dat associated with the *FTrees* installation. For example:  
C:\Programs\BioSolveIT\FTrees2\config\_ft.dat

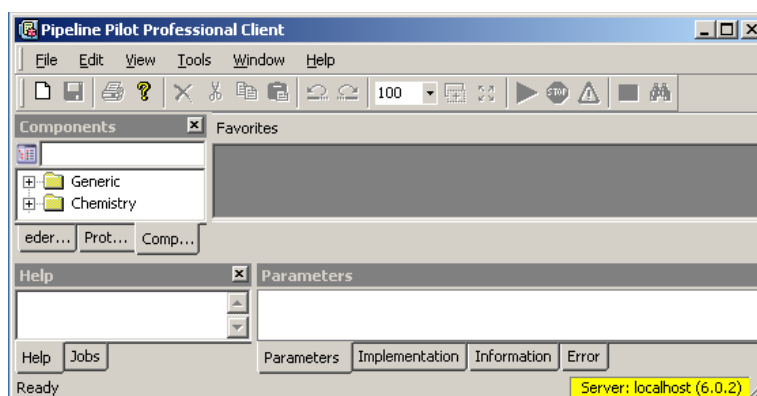


Figure 3.2: See where your Pipeline Pilot Server is installed.

You can save these settings in the components.

### 3.2.2 To Connect to an *FTrees* Installation on a Remote Linux Server

Your existing *FTrees* installation could be on a Linux computer remote from the Pipeline Pilot server – in this case we offer an alternative so you can use the remote installation instead. Here, the calculations carried out by *FTrees* will be done on the remote Linux machine. Pipeline Pilot logs into the Linux machine using `ssh`, having copied all relevant files to the machine, and will run the calculation there, finally copying back all data it needs to the Pipeline Pilot server to continue with the pipeline.

This method is given the name *via SSH*. You can see an example in figure 3.3.

Implementation	
User-Mode	Default
<input checked="" type="checkbox"/> Run <i>FTrees</i>	via SSH
<input checked="" type="checkbox"/> via SSH	
Executable (ft)	/software/BioSolveIT/ftrees/ftrees
Configuration	/home/ederk22s/config_ft.dat
Host	rho
User	ederk22s
Password	*****

Figure 3.3: *via SSH*: Connect to an *FTrees* Installation on a Remote Linux Server via `ssh`

- Requirement:
  - *FTrees* is installed on a Linux machine available to the Pipeline Pilot server via `ssh`.
- General Steps:
  1. Set <Run *FTrees*> to *via SSH*

2. For the parameter <Run FTrees-> via SSH -> Executable>, enter the path to the *FTrees* installation on the Linux machine. For example:  
/software/BioSolveIT/ftrees-2.0.1/bin/ftrees-2.0.1
  3. For the parameter <Run FTrees-> via SSH -> Configuration>, enter the path to the *FTrees* configuration file config\_ft.dat associated with the *FTrees* installation. For example:  
/software/BioSolveIT/ftrees-2.0.1/config\_ft.dat
  4. For the parameter <Run FTrees-> via SSH -> Host>, enter the Linux machine host name
- User specific steps:
    1. For the parameter <Run FTrees-> via SSH -> User>, enter the user login name for ssh on the Linux machine
    2. For the parameter <Run FTrees-> via SSH -> Password>, enter the user password for ssh on the Linux machine
    3. There are more advanced options to be found under <Run FTrees-> via SSH -> Options> for more specific ssh parameters. Note the option <Run FTrees-> via SSH -> Options -> Delete Results> may be useful for trouble-shooting later.

You can save these settings in the components (be sure not to save your own user specific login details in components available to others!).

All files necessary for the *FTrees* calculation will be transferred via `scp` between the Pipeline Pilot server and ssh Linux machine. Files copied and files created the remote server are automatically deleted at the end of the job leaving no trace. However, in case the user would like to leave a copy of the calculation and result files on the Linux machine, or for trouble-shooting as mentioned above, it is possible to set a parameter to tell Pipeline Pilot not to delete these files:

<Optionally> : <sshDeleteResults> : *False*

# Trouble Shooting

## 4.1 Problems connecting to the external installation

Other most commonly seen problems with *FTrees* in Pipeline Pilot is with the connection to the external *FTrees* installation. For one thing, *FTrees* itself must be correctly installed on the system independently of Pipeline Pilot – it is essential first to make sure this is the case (especially to make sure that *FTrees* can locate the licenses). Once *FTrees* runs fine on your system, the remaining key task is to make sure the paths to the executable and configuration file are correct within the Pipeline Pilot components.

If something is amiss with the connection to the external installation, you will see an error message box pop up. Check the paths to the executable or to the configuration file. If it seems *FTrees* could be started but not run, the problem almost always lies either with the path to the configuration file or with licenses not being found. **Also make sure that you have a *FTrees* license [FFS].** If *FTrees* runs OK independently from Pipeline Pilot then it is likely to be the path to the configuration.

When the error messages pop up, they may contain an `FTreesError` in the error message box, as in figure 4.1.

Go to the 'Jobs' tab below the Protocol workspace and check under the last run job for a file called 'FTreesComponent Debug' – as in figure 4.2. Clicking on the link brings up HTML report with input and output data in a browser.

A correctly started *FTrees* job outputs the following header: if there is a problem you will see some of this header and the point where the problem occurs:

```

                                     F T R E E S
                                     Feature-based molecular similarity finder
=====
Version:  2.0.2  (11.08.08)
Modules:  [PVM] [FFS]
/ written by: Matthias Rarey, Marc Zimmermann,
              Sally Hindle, Robert Fischer
```

```

/
/ copyright by: BioSolve IT GmbH, FhG -SCAI
/               Sankt Augustin, Germany
/ for further information mail ftrees-info@biosolveit.de /
/
Additional copyright notes:
  getline library: (C) 1993 by Chris Thewalt
  PVM library: (C) 1997 by University of Tennessee, Knoxville TN
>> FTrees configuration file '/software/BioSolveIT/FTrees/ftrees2.0.1/
  config_ft.dat' loaded.

>> Licensed modules: FTrees [PVM] [FFS]
>> PVM status: no pvm daemon, running sequential.
>> Scripts are executed in sequential mode; start PVM for parallel mode.
>> SETTINGS = '/software/BioSolveIT/ftrees2.0.1/static_data/ftrees_settings.dat'
  loaded.
>> CHEMPAR = '/software/BioSolveIT/ftrees2.0.1/static_data/chempar.dat' loaded.
>> CONTACT = '/software/BioSolveIT/ftrees2.0.1/static_data/contact_ft.dat' loaded.
>> TRANSFORM = '/software/BioSolveIT/ftrees2.0.1/static_data/transform.dat' loaded.
>> GRAPHIC = '/software/BioSolveIT/ftrees2.0.1/static_data/graphic_ft.dat' loaded.
...

```

## 4.2 Problems using the *ssh* Method

You may also experience problems using the *ssh* login, for example, the user name is unknown or the host is not found.

## 4.3 Further help

More complicated errors may arise during the running of *FTrees*. Again though, the errors will be collected and as much information shown as possible. If you are familiar with *FTrees* you may want to take a look at all the output of the job yourself to see if you can recognize the problem. In this case, you can look in the temporary folders Pipeline Pilot sets up internally to find the output, or if you are working with the *ssh* method, set the parameter <Run FTrees -> via SSH -> Options -> Delete Results> to *False* so that you may then find the files retained on the *ssh* host: these will be in the directory set under the *ssh* parameter <Run FTrees -> via SSH -> Options -> Temp Path> (see the help text associated with this parameter to find its default value) – essentially a cryptically named folder whose name begins with the date and time of the job!

If you still do not know what is causing the errors, write down as much information as possible relating to your installation scheme and cut and paste the text from the 'Errors.txt' file or, if possible, the output files from the *FTrees* job itself into an email. Send all the information to:

support@biosolveit.de

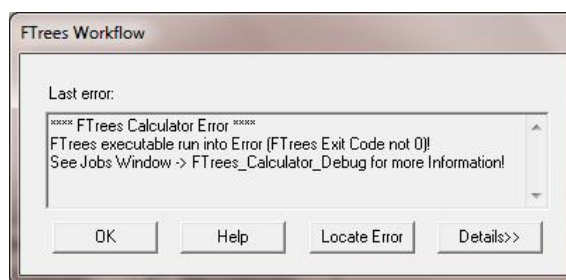





Figure 4.1: An error box reporting that the *FTrees* exe could not be found.

Jobs	
Protocol Name	Status
New Protocol1	Finished
 FTrees Workflow	Error
├─ <a href="#">Errors.txt</a>	
└─ <a href="#">FTrees_Calculator_Debug...</a>	
 New Protocol1	Running
 New Protocol12	Finished

Jobs [Help](#) [Error](#)

Figure 4.2: The full error report can be found in the 'Jobs' tab.



# Tips and Tricks

## 5.1 Other Significant Parameters in the *FTrees* Components

For detailed documentation of all parameters for all components, refer to the documentation you find in the *Help* area of the Pipeline Pilot window.

We list here particularly interesting parameters: those that greatly influence the workflow or change the outcome of calculations, or those that may help you understand what is happening in the components.

### 5.1.1 <Has Same File System>

*FTrees*Calculator/Similarity : Implementation : <Run *FTrees*-> via SSH -> Options -> Has Same File System>

The name of this parameter stands for "ssh Has Same File System". Normally, for an *ssh* job, Pipeline Pilot must first copy all the data required by *FTrees* to the <Run *FTrees* -> via SSH -> Host> using *scp*. This is time consuming. It is possible that the Pipeline Pilot Server and <Run *FTrees* -> via SSH -> Host> actually share the same file system rendering the *scp* process unnecessary. Select True if the Pipeline Pilot Server and <Run *FTrees* -> via SSH -> Host> share the same File System - no copying of data is necessary. Selecting False means Pipeline Pilot copies all the data to and back from the <Run *FTrees* -> via SSH -> Host> - this is just a little slower but will always still work. Leave the parameter set to False if you are uncertain!

## 5.2 Accessing Other Domains within Pipeline Pilot

Often in house data or even your own working data are accessible from a windows computer via a domain (a path starting for example 'z:\...' or '\\...\') which you cannot find from within Pipeline Pilot. That means you must first literally transfer the data to the Pipeline Pilot Server itself.

If you are using a Linux Pipeline Pilot server, this hint does not apply.

To get around this problem and make the Pipeline Pilot working environment much more flexible you can allow users access to domains – you need Administrator rights to be able to do this! Also check first that you should change these settings as they may have already been set to fit the current environment.

- Go to the 'Scitegic Server Home Page', for example, via the Help menu in your Pipeline Pilot client.
- Click on 'Pipeline Pilot Administration Portal' and log in with the Administrator user name and password.
- In the Security tab go to Authentication.
- For the 'Authentication Method' choose 'DOMAIN' and a set of parameters will appear.
- Enter the domain name in the field 'Domain' and choose 'Full' for 'Impersonation'
- Choose 'DOMAIN' for 'Retrieve Groups' and leave 'Limit access to listed domains' set to 'No'
- click 'Save' and log out again.

After you have done this you will need to enter your domain login details when you start the Pipeline Pilot Client.

# Bibliography

- [1] Judd D. B. Watson S. P. Lewell, X. Q. and M. M. Hann. Recap - retrosynthetic combinatorial analysis procedure: A powerful new technique for indentifying privileged molecular fragments with useful applications in combinatorial chemistry. *Journal of Chemical Information and Computer Science*, 38:511–522, 1998. 9
- [2] M. Rarey and J.S. Dixon. Feature trees: A new molecular similarity measure based on tree matching. *Journal of Computer-Aided Molecular Design*, 12:471–490, 1998. 7
- [3] M. Rarey and M. Stahl. Similarity searching in large combinatorial chemistry spaces. *Journal of Computer-Aided Molecular Design*, 15:497–520, 2001. 2, 7